

ABSTRACT

As the pace of novel fold discovery has slowed, and as structure determination has become more routine, there is a growing degree of redundancy within the PDB, whereby a given fold or protein is represented by multiple structures. Coincident with these trends, the rapidly growing volumes of data being produced by next-generation sequencing initiatives provide both questions and opportunities relevant to protein conservation at the residue level: regions within proteins may exhibit human-specific or cross-species conservation, but for which no clear structural mechanism is acting as the evolutionary constraint. Given the complex nature, fundamental importance, and general ubiquity of allosteric regulation (Gunasekaran et al, 2004; Tsai et al, 2014), allostery provides a rational framework as a starting point to understanding these otherwise cryptic sites.

Herein, we describe a workflow to automate the identification of structures that occupy alternative energetic minima, and apply this framework to the wealth of data available in the PDB. We aim to identify those key residues that may mediate allostery or otherwise function as significant regulatory sites in the context of conformational change. We describe a refinement of the binding leverage concept introduced by Mitternacht and Berezovsky, to identify essential patches on the surface, as well as a network-based formalism employing models of conformational transitions to identify essential residues which can function as information flow bottlenecks, which may in turn link the surface patches. Notably, this framework is highly mechanistic in nature, in contrast to approaches that use conservation or biophysical properties of the individual residues. We demonstrate that the sites identified are generally conserved (both across species and within human populations), and identify cases in which these sites coincide with human disease loci. Finally, we provide a user-friendly interface to perform these predictions on submitted protein structures.

INTRODUCTION

In addition to the conformations exhibited by an ensemble of structures, the underlying energetic landscape of a protein is dynamic in nature: allosteric signals or other external changes may actually reconfigure and reshape the landscape, making

- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: TITLE NEEDED – suggested by KKY: Identification of allosteric residues by exploring alternative conformations in PBD ... [1]
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: Throughout the early years of protein structure determination efforts, the primary objective has been the discovery of novel folds and structures. However, as
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: easier
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: This redundancy has partially enabled structural genomics initiatives, and has also shifted the focus from novel fold discovery to understanding protein dyn ... [2]
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: there is a
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: appreciation of the role that ... [3]
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: play in virtually all
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted:).
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: Using the proteins identified ... [4]
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: the
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: might
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: One approach is
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: et al
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: and the second is
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: method
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: internal
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: , some of
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: these complementary methods are
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: relying on
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: Early models of protein stru ... [5]
- DECLAN CLARKE 7/19/15 12:42 AM
Deleted: itself

previously high-energy configurations more favorable, for instance, thereby shifting the relative populations of states within an ensemble (Tsai et al, 1999). Thus, energy landscape theory provides the conceptual underpinnings necessary to describe how proteins change behavior under changing conditions and stimuli.

The growing appreciation for dynamic behavior and the importance of conformational heterogeneity is being facilitated by a saturation in the number of folds (Supp. Fig. 1a) and a growing redundancy within the PDB. Such redundancy is represented, for instance, when the same protein is structurally resolved under different conditions, potentially resulting in alternative conformations. This redundancy has also aided structural genomics initiatives (for cases in which template-based modeling is used to predict the structure of a protein that would otherwise be difficult to crystallize; Ashkenazy et al, 2011), in addition to motivating efforts to catalogue and model conformational changes on a database-wide scale (Gerstein et al, 1998; Krebs et al, 2000; Echols et al, 2003; Flores et al, 2006).

Given that a primary driving force behind shaping the evolution of these landscapes is the need to form highly efficient macromolecules that can regulate activity according to environmental signals, the mechanisms responsible for allostery and conformational change in general have been given greater attention. An allosteric mechanism may involve the modulation of large-scale motions upon binding of an effector ligand, with this modulation resulting in changes in conformation or dynamic behavior of a ligand-binding that is typically very distant from the effector-binding site. Such motions may also affect patterns of communication between residues internal to the protein, thereby enhancing, modulating, or rewiring essential channels of information flow within the protein.

Interestingly, deep sequencing has unearthed a class of conserved residues for which no clear structural constraints seem to be responsible for such conservation (i.e., “cryptic sites”). Given the complex nature and difficulty associated with studying allosteric regulation, we note that the results outlined in our study may help to explain the conservation of such residues. We also find that our identified sites tend to be conserved across species.

PLAT

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: also

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: in

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: In addition to being facilitated by newer experimental methods, a need to describe allosteric behavior, and the introduction of energy landscape theory, this

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: trend toward sequence

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: This growing redundancy shifts the emphasis from static macromolecular architectures to structural heterogeneity and ensembles – *a priori*, X-ray crystal structures constitute proteins in energetic wells, and may thus provide snapshots of alternative conformations that, in turn, exhibit different functional characteristics (such as different catalytic rates or ligand specificities). As the volume of crystal structures continues to expand, there is a growing need and opportunity to leverage this data to identify alternative conformations, and to thus more comprehensively elucidate their energetic landscapes and modes of regulation.

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: impart allosteric regulation

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: within the

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: surface site

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: It

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: entail changes in how different sub-regions within a

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: interact, with some allosteric signals

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: mutual interactions between regions, and other signals dampening such interactions

Numerous methods have been devised for the prediction of allosteric residues. Conservation has been used as a principle metric in the search for allosteric sites, either in the context of conserved residues (Panjkovich and Daura, 2012), networks of co-evolving residues (Lee et al, 2008; Suel et al, 2003; Lockless and Ranganathan, 1999; Shulman et al, 2004; Reynolds et al, 2011; Halabi et al, 2009), or conservation in structure (Panjkovich and Daura, 2010). In related studies, both conservation and geometric-based searches for allosteric sites have been successfully applied to a few systems (Capra et al, 2009), several of which also employ SVMs (Huang et al, 2006, Huang et al, 2013). Panjkovich and Daura have also used normal modes analysis, coupled with ligands of varying size, to examine the extent to which bound ligands interfere with low-frequency motions (Panjkovich and Daura, 2012). Similar approaches for identifying surface sites that may be allosteric have also been explained by others (Mitternacht and Berezovsky, 2011; Ming and Wall, 2005).

The importance of protein dynamics to allosteric function has been exploited for identifying allosteric sites in several other studies. Specifically, normal modes have also been used extensively by the Bahar group in order to identify important subunits of proteins that act in a coherent manner for specific proteins (Chennubhotla and Bahar, 2006; Yang and Bahar, 2005). Rodgers et al applied normal modes to identify and experimentally validate the importance of key residues in the CRP/FNR transcription factors (Rodgers, 2013). Several groups have applied molecular dynamics and communities-based network analysis in order to identify internal residues which may function as allosteric bottlenecks (Sethi et al, 2009; Gasper et al, 2012; VanWart et al, 2012; see also reviews by Csermely et al, 2013, as well as Rousseau and Schymkowitz, 2005). In conjunction with NMR, Rivalta et al use molecular dynamics with the same communities-based network analysis to identify essential elements in imidazole glycerol phosphate synthase (Rivalta et al, 2012).

Though having provided valuable insights, many of these approaches may be limited in terms of application (i.e., with respect to the subclass of protein studied), scale (the numbers of proteins which may be feasibly investigated), or the class of residues identified (i.e., surface or interior).

With the objective of focusing on those proteins that may undergo conformational change as part of their allosteric behavior, we develop and apply a general framework to identify instances of alternative conformations, and apply this method to the entire PDB. We then determine sites within the protein that could potentially affect the thermodynamic stability of these conformational states in a computationally tractable manner. Specifically, allosteric regulation is usually imparted by a small subset of residues within a given protein: these may constitute a set of residues on the protein surface (such as an effector binding site) or a channel of residues within the interior that are responsible for linking distal sites (i.e., bottleneck residues). Our objective is to identify both of these classes of residues within a unified study. Several of our identified sites correspond to human disease loci for which no clear mechanism had previously been proposed.

Several of the advantages of our method include the fact that it is mechanistic in nature (as opposed to using less direct measures such as conservation), it may be applied to a wide variety of proteins and on a mass scale, and it simultaneously captures both surface sites and many of the interior residues which may bridge such sites. Finally, our pipeline (termed STRESS, for STRucturally-identified ESSential residues) is made available to the public through a server hosted by our group, from which users may submit their own structures for analysis.

RESULTS

High-Throughput Identification of Structures in Distinct Energetic Wells

To build the high-confidence dataset of conformational changes, we use a generalized approach to leverage the wealth of data in the PDB for systematically identifying proteins that occupy alternative energetic wells.

As a first step toward culling a high-confidence set of alternative conformations, we perform multiple structure alignments (MSAs) across sequence-identical domains as well as proteins, with these structures having been filtered by resolution and other metrics to ensure quality (see Methods and Fig. 1 for details). We first worked with domains to probe for intra-domain conformational changes of functional significance. In addition, better structure alignments are generally possible at the domain level. The filtered dataset

WORD 2.

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: The elucidation of allosteric behavior in the context of such conformational shifts has traditionally been obtained through experiments on proteins that are limited in number, size, or both. Such studies [[cite Ranganathan, maybe also MD studies, others]] achieve high accuracy and yield valuable insights, but the limitations in scope and applicability of in-vitro studies on individual systems have made the systematic study across thousands of proteins infeasible. - ... [6]

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: A mechanistic approach is applied for identifying surface residues that may be allosterically significant, and we apply a dynamical networks-based community analysis to identify residues that may be essential for transmitting allosteric information, most of which are internal to the protein. The sites identified are generally conserved, both across species and even within human populations.

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: through

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: It may be that a subset of these proteins do not exhibit allosteric behavior as part of their native functionality within cells, but the multiple energetic minima captured in their crystal structures may nevertheless be exploited for protein engineering [[cite C.J. Wilson, others]] or in pharmaceutical contexts [[cite]] (see preceding discussion).

of domains contains 79% of all available crystal structures in the PDB (as of December 2013). PDB-wide MSAs across sequence-similar groups reveal that, in agreement with expectation, average pairwise root-mean-square deviation (RMSD) values increase at lower levels of sequence identity, as do Q_H values (Q_H , an alternative metric to RMSD, quantifies the degree to which residue-residue distances differ between two conformations, and is detailed in [cite] and Methods) (Supp. Fig. 2).

After performing multiple structure alignments for each sequence-identical group of proteins, we use the resultant pairwise RMSD values to infer distinct conformational states. This is a non-trivial task, as considerable sources and degrees of variation are inherent in the MSA of each protein (sources of such variation may lie in either experimental limitations or the biology itself; examples include limited resolution, the fact that small differences in ensembles of structures may occur within a single energetic well, or slight differences in crystallization conditions). Thus, in order to reduce false positives and confidently automate the identification of biologically relevant, large-scale, and truly distinct energetic wells, we apply a modified K -means clustering algorithm in order to assign each structure to a particular well within an MSA; structures that cluster together (i.e., exhibit low pairwise RMSD) constitute a given well (Fig. 2B; see Methods).

This algorithm (termed K -means clustering with the gap statistic) identifies the ideal number of clusters (i.e., K) to describe a dataset in an automated way by comparison with a randomized null dataset. The algorithm is further detailed in methods and in (Tibshirani, 2001). Briefly, the K -values obtained using this algorithm, which we take to represent the number of distinct energetic wells, is used to reduce the uncertainty associated with limited crystallographic resolution and the potential for a protein to exhibit subtle conformational heterogeneity within a single well.

Modified Binding Leverage to Identify Critical Residues on the Surface

We use the high-confidence set of alternative conformations described above as the input for predicting the residues that are most important in allosteric regulation and activity. In the first, **effector** binding sites on the protein surface (some of which may act as latent ligand binding sites **and active sites**) are identified using a modified version of

DECLAN CLARKE 7/19/15 12:42 AM

Formatted: Subscript

DECLAN CLARKE 7/19/15 12:42 AM

Formatted: Subscript

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: k

DECLAN CLARKE 7/19/15 12:42 AM

Deleted:

DECLAN CLARKE 7/19/15 12:42 AM

Deleted:

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: To validate the output generated by this clustering algorithm, we manually annotated the MSAs of several well-studied canonical allosteric systems. The gap statistic performed well in discriminating crystal structures that were manually determined to constitute alternative biological states. Some of the key systems we manually annotated include proteins that have been very well-studied and characterized in the literature, such as tyrosine phosphatase (Wiesmann et al, 2004), DNA polymerase I (Xiang et al, 2006), adenylate kinase (Arora et al, 2007), Hsp ATPase (Liu et al, 2010), phosphoglycerate dehydrogenase (Grant et al, 1996), phosphofructokinase (Laurent et al, 1984), phosphotransferase (Kohl et al, 2005), and alanyl-tRNA synthetase (Dignam et al, 2011). For each of these cases, we manually determined that there were two main biological states (for example, with and without bound ligand, Supp. Fig 4). The gap statistic correctly determined that the appropriate K value for these cases was two. - ... [7]

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: to two complementary methods

the “binding leverage” framework for ligand binding site prediction originally described in 2011 (Mitternacht and Berezovsky, 2011, also detailed in Methods).

Briefly, this method first entails using a series of Monte Carlo simulations to probe the protein surface (with the protein being represented with all heavy atoms) with a simulated ligand, thereby generating a series of candidate sites. Each candidate site is then scored on the basis of the degree to which the occlusion (with the simulated ligand) disrupts the large-scale motions of the protein (Fig. 1, bottom left; see Methods). These motions are taken from anisotropic network models, but the results do not change drastically if we directly use the alternative conformations as given by the crystal structures (see “Comparisons Between Different Models of Protein Motions” in Results).

The main modifications to the formalism described by Mitternacht and Berezovsky include the inclusion of heavy atoms in the protein during the Monte Carlo search, in addition to an automated means of thresholding the list of ranked sites to give a more selective set of candidate sites. This approach results in finding an average of ~2 distinct binding sites per domain (Fig. 2a; see Methods for the details on defining distinct sites).

Surface residues important to allosteric behavior may either be the binding sites for allosteric ligands or the allosteric ‘sinks’ in signal transmission (that is, the sites affected by binding to a ligand in a distal region). In order to evaluate the extent to which this method identifies sites of the former category, we studied the ligand-binding sites and active sites in a set of 12 well-studied systems for which the crystal structures of both the *holo* and *apo* states are available (Supp. Table 1; note that these 12 systems differ from the proteins highlighted in our alternative conformations analysis, as that set of proteins highlights our ability to identify alternative conformations in many different biological contexts, and not just ligand binding – see "Identifying Distinct Conformations in an Ensemble of Structures" in Methods and also Supp. Fig. 4). In order to be consistent with the original binding leverage study, these 12 constitute those used by Mitternacht and Berezovsky for which the *apo* and *holo* states are clear (Supp. Table 1).

We find that, out of the 12 canonical systems studied, we positively identify an average of 60% of the sites known to be directly involved in ligand or substrate binding. It has previously been shown that it is especially difficult to identify the sites in aspartate transcarbamoylase (Mitternacht and Berezovsky, 2011); excluding aspartate

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: developed by

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: et al (Mitternacht et al)

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: an artificial

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: artificial

DECLAN CLARKE 7/19/15 12:42 AM

Deleted:)

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: et al

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: Supp.

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: 5a

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: Table 1).

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: -

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: sites

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: et al

transcarbamoylase from this analysis results in finding an average of 65% of known biological sites. We note that these statistics are achieved by covering an average of 15% of proteins' residues (Supp. Table 2). For most proteins, selecting 15% of the residues is conservative -- more than 15% of the proteins' residues are involved in ligand or substrate binding for most proteins (Supp. Table 3).

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: ligand-binding

Some of the sites identified do not meet the thresholds needed for defining a site of known biological significance. However, two factors suggest that such sites may still be significant in an allosteric context. First, we emphasize that such sites may nevertheless correspond to latent allosteric regions (Bowman et al, 2015): even if no known biological function is assigned to such sites, their occlusion may still disrupt large-scale motions. Such latent allosteric pockets may be useful in the context of drug development and targeting. Secondly, we often find that these sites nevertheless exhibit some degree of overlap with sites of biological interest, suggesting that the identified sites often lie within the neighborhood of known biological sites (Supp. Table 4).

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: We note that, for those

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: that constitute false positives (sites which we predict to be important for allostery, but which nevertheless

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: ligand binding sites), we often find overlap with sites of biological interest, suggesting that the identified sites often lie within the neighborhood of known biological sites (Supp. Table 4). We also

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: our high-scoring sites which do not correspond to known biological ligand-binding

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: sites

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: regions, the

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: of such sites

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: (Here,

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: ,

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: ,

DECLAN CLARKE 7/19/15 12:42 AM
Deleted:).

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: ConSurf

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: BL-

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: ConSurf

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: . The significance of the disparity was

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: (using a

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: test

Conservation of Surface Sites Across Species

Residues that lie in our prioritized sites tend to be more conserved, on average, than other residues of the same protein with the same degree of burial (Fig. 3C). Both the In order to be consistent with our networks-based analysis for identifying essential interior bottleneck residues (termed "interior critical residues" throughout), the degree of each residue (representing the number of other residues with which that residue interacts) is used to characterize burial. The conservation is evaluated using ConSurf scores (termed "conservation scores" throughout; Glaser et al, 2003; Landau et al, 2005; Ashkenazy et al, 2010; Celniker et al 2013), and these results constitute the distribution of scores for proteins in our entire dataset. Here, surface critical residues had an average conservation score of -0.131, whereas non-critical residues with the same degree distribution (i.e., same degree of burial within the protein) had an average score of +0.059, demonstrating that critical residues tend to be more conserved ($p < 2.2e-16$, Wilcoxon rank sum).

1000 Genomes- and ExAC-Derived Conservation of Surface Sites Amongst Modern-Day Humans

EMBARGO

Although we observe a general trend in which rare alleles coincide with surface critical residues, the trend is not significant at the level of 0.5% (Fig. 3B). The lack of significance may partly be a consequence of the limited number of proteins (44) in our dataset that are hit by 1000 Genomes single-nucleotide variants (SNVs). In addition, we note that the long tail extending to lower allele frequencies for critical residues may suggest the possibility that only a subset of residues in our prioritized binding sites are essential (i.e., each one of our sites has 10 residues, but there may only be a small subset of these 10 which are important allosterically, thus explaining the long tail toward lower DAF values in Fig. 3B). As an alternative test, we also performed an analysis of the potential shifts in these distributions using a two-sample Kolmogorov-Smirnov test ($p=0.08$, Supp. Fig. 13a).

We also performed an analogous analysis using the data provided by the Exome Aggregation Consortium (Exome Aggregation Consortium, abbreviated ExAC). The results obtained using ExAC (in which we use AF values instead of DAF – see Methods) are similar to those using 1000 Genomes data (Supp. Fig. 10). Although the mean minor allele frequencies for surface-critical residues are higher than those of non-critical residues (Supp. Fig. 10, left), there is a skew toward lower minor allele frequencies for critical residues relative to non-critical residues. We also point out that the median minor allele frequency for surface-critical residues is substantially lower than that for non-critical residues. The relative shifts of these distributions are also shown in Supp. Fig. 14a ($p=0.0475$, two-sample Kolmogorov-Smirnov test).

In addition to examining the distribution of derived allele frequencies, we also examined the *fraction* of rare alleles (defined as the ratio of the number of low-DAF SNVs to all SNVs in a given protein for 1000 Genomes, or alleles with low minor allele frequency in ExC) in critical residues and non-critical residues for those proteins for which at least 1 critical residue and 1 non-critical residue is hit by an SNV. Using different DAF cutoffs for 1000 Genomes variants (0.5% and 0.1%) to define rarity, the results for surface residues are summarized in Supp. Fig. 7. Using a rarity threshold of 0.5% (0.1%), the fraction of rare SNVs in critical residues exceeds that in non-critical residues for 6 (16) structures (in green), whereas the fraction of rare alleles in non-critical surface residues exceeds that in critical residues for only 2 (8) structures (in gray).

TOO LONG

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: 05

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: SNPs.

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: be indicative of

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: 3B

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [1]

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [2]

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [3]

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: SNPs

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: SNPs

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: a 1000 Genomes non-synonymous SNP.

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: (0.05

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: 01

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: 05

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: 01

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: SNPs

The fraction of rare alleles with ExAC provided similar results; we find that surface residues are generally more conserved than other residues, and this result holds for different thresholds for defining rarity (Supp. Table 7). For surface-critical residues, 9.5% (30.0%) of the structures studied were such that the fraction of rare alleles in critical residues exceeded that for non-critical residues using a threshold of 0.005 (0.001), whereas the opposite trend was observed in only 6.0% (13.0%) of eligible structures.

Finally, we note that SIFT and PolyPhen scores of critical and non-critical residues hit by variants from the ExAC dataset were also evaluated, as this may provide information regarding the effects of SNVs which is orthogonal to examining minor allele frequencies alone. No significant disparity was observed between critical and non-critical residues when evaluating SIFT scores (Supp. Fig. 11, left). However, surface critical residues were shown to exhibit significantly higher PolyPhen scores relative to non-critical residues, thereby demonstrating the more deleterious nature of SNVs which fall on critical residues (Supp. 12, left; note that higher PolyPhen scores denote more damaging variants).

Identifying Internal Residues for Transmitting Allosteric Information Using Dynamical Network Analysis

The binding leverage framework described above captures hotspot regions close to or at the surface of the protein, but the Monte Carlo search employed is *a priori* excluded from the protein interior. Thus, motivated by previous studies focused on individual proteins, such as tRNA synthetase (Sethi et al, 2009), essential metabolic enzymes (Manley et al, 2013), and the HIV envelope glycoprotein (Sethi et al, 2013), we apply communities-based network analyses to the protein complexes of our dataset to identify important internal residues, in addition to residues that may be closer to the surface, and note that the identification of residues in the protein interior are often needed to serve as communication channels between surface sites.

The nodes of the network represent individual residues, and edges between these nodes are drawn between residues that lie within a mutual proximity of 4.5 Angstroms. The first step in this analysis is to define the communities within the network. A given “community” refers to a group of residues that are highly inter-connected, but which have

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: ExAC-Derived Conservation of Surface Sites Amongst Modern-Day H... [8]

DECLAN CLARKE 7/19/15 12:42 AM

Moved up [1]: results obtained using ExAC (in which we use AF values instead of DAF – see Methods) are similar to those using 1000 Genomes data (Supp. Fig. 10).

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: Although the mean allele frequencies for surface-critical (BL) residues are higher than those of non-critical residues (Supp. Fig. 10, left), there is a skew toward lower allele frequencies for critical residues relative to non-critical residues. We also point out that the median allele frequency for surface-critical residues is substantially lower than that for non-critical residues. ... [9]

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: .

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: .

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: (Supp.

DECLAN CLARKE 7/19/15 12:42 AM

Moved up [2]: Fig.

DECLAN CLARKE 7/19/15 12:42 AM

Moved down [4]: Fig.

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: bridge the gap

minimal edges to residues outside the community (see Methods). We have tried applying both an information-theory based method (termed “Infomap”, Rosvall et al, 2007) and the classical Girvan-Newman (GN) formalism (Girvan et al, 2002) to decompose networks of interacting protein residues into communities (see Methods), and find that GN is more appropriate (see Supplemental discussion). The results described here are thus based on the GN formalism.

METH

In order to recapitulate the contributions of various edges to information flow over the course of conformational changes, edges are weighted on the basis of the correlated movements using anisotropic network models (see [Methods and Yang et al, 2005](#)); if two contacting residues exhibit a high degree of correlation, then this implies that the motion of one residue may tell us about the motion of the other residue. This suggests a strong flow of energy or information between the two residues, and when calculating the path distance between the two residues, a short distance is assigned (this lowers the shortest paths in which this pair of residues participate, thereby favorably weighting this pair as a potential bottleneck within the network— see Methods for more details). Finally, once all connections between contacting pairs are appropriately weighted and the communities are assigned, a residue is deemed to be critical for allosteric signal transmission if it is involved in a highest-betweenness edge connecting two distinct communities (see Methods). For instance, applying this method to threonine synthase results in the community partition and associated critical residues highlighted in Supp. Fig. 6.

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: :

Conservation of *Interior Critical Residues Across Species*

Residues which function as essential allosteric conduits of information in mediating signal transduction from one site of a protein to another are likely to be more conserved, on average, than residues which lie outside of such channels. Thus, as for the case with [surface](#) critical residues, [we](#) evaluate conservation of residues identified as critical by GN, and compare these conservation scores to those of non-critical residues with the same degree (Fig. 3F).

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: GN

Again, each point in this plot represents the mean [conservation](#) score for all residues within one of two classes (the set of [interior](#) critical residues or the randomly-selected set of non-critical residues with the same degree) within a particular protein

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: BL

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: as a validation of our method, we use the ConSurf server to

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: ConSurf

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: GN

structure. The randomly-selected non-critical set of residues was chosen in a way such that, for each critical residue with degree K, a randomly-chosen non-critical residue with the same degree K was included in the set. The distribution of non-critical residues shown is very much representative of the distribution observed when re-building the random set many times.

Interior critical residues are generally found to be more conserved than non-critical residues with the same degree of burial: the average **conservation** score for **interior** critical residues is -0.179, whereas that for non-critical residues with the same degree is -0.102 ($p=3.67e-11$, Wilcoxon rank sum test).

1000 Genomes- and ExAC-Derived Conservation of **Interior** Residues Amongst Modern-Day Humans

The residues we identify as **bottlenecks** are shown to be under negative selection in the context of modern-day humans. Our analyses of variants identified from The 1000 Genomes shows that **interior** critical residues occur at sites with significantly lower DAF values (Fig. 3E). Similar results are obtained for ExAC: **interior** critical residues exhibit a significantly lower minor allele frequency than do non-critical residues (Supp. Fig. 10, right). We also performed an analysis of the potential shifts the distributions of the mean DAF values using a two-sample Kolmogorov-Smirnov test for 1000 Genomes ($p=8.9E-5$, Supp. Fig. 13b) and ExAC variants (Supp. Fig. 14b, $p=8.7E-5$).

We examined the fraction of rare alleles in critical **interior** residues and non-critical residues for those proteins for which at least 1 critical residue and 1 non-critical residue is hit by an SNV. The corresponding results for 1000 Genomes and ExAC variants are given in Supp. Figs. 8 and Supp. Table 7, respectively. Using a rarity threshold of 0.5% (0.1%) with 1000 Genomes variants, the fraction of rare SNVs in critical residues exceeds that in non-critical residues for 10 (20) structures (in green), whereas the fraction of rare allele in non-critical surface residues exceeds that in critical residues for only 0 (1) structure (in gray) Fig. When measuring human-specific conservation using the fraction of rare alleles in ExAC, interior-critical residues are also more conserved than non-critical residues at varying thresholds for rarity (Supp. Table 7). 15.5% (41.1%) of the structures studied were such that the fraction of rare allele in

COMBINE WITH CONSERVATION

+

- DECLAN CLARKE 7/19/15 12:42 AM Deleted: GN
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: ConSurf
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: GN
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: GN
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: critical
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: GN-
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: 3E). ... [10]
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: , we
- DECLAN CLARKE 7/19/15 12:42 AM Moved (insertion) [5]
- DECLAN CLARKE 7/19/15 12:42 AM Formatted: Indent: First line: 0.5"
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: analyzed
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: network
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: a 1000 Genomes non-synonymous SNP.
- DECLAN CLARKE 7/19/15 12:42 AM Moved up [3]: Fig.
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: 8.
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: 05
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: 01%),
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: SNPs
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: ... [11]
- DECLAN CLARKE 7/19/15 12:42 AM Moved (insertion) [4]
- DECLAN CLARKE 7/19/15 12:42 AM Deleted:
- DECLAN CLARKE 7/19/15 12:42 AM Moved up [5]: 10, right).
- DECLAN CLARKE 7/19/15 12:42 AM Deleted: ... [12]

critical residues exceeded that for non-critical residues using a threshold of 0.005 (0.001), whereas the opposite trend was observed in only 0% (3.3%) of structures.

SIFT and PolyPhen scores of interior critical residues hit by variants from the ExAC dataset exhibited trends similar to those observed for surface-critical residues: no significant difference was seen between interior critical and non-critical residues with respect to SIFT scores (Supp. Fig. 11, right). However, these critical residues were shown to exhibit significantly higher PolyPhen scores relative to non-critical residues (Supp. Fig. 12, right), suggesting that modifications to these critical residues are significantly more damaging.

Comparisons Between Different Models of Protein Motions

Given that our entire scheme is based on an understanding of protein motions, we evaluated the extent to which the results may be sensitive to different models of conformational change. ANMs are simple and straightforward to apply on a database scale, and are thus used as our primary model of choice.

As an alternative to ANMs, one may simply use the displacement vectors between all corresponding pairs of residues within the two crystal structures of the alternative conformations for a given protein. This more direct model of conformational change, which we term absolute conformational transitions (ACT), may be applied in a straightforward manner to single-chain proteins. When we use ACT to apply the modified binding leverage framework for such single-chain proteins, we observe that our surface critical residues are significantly more conserved than are non-critical residues (Supp. Fig. 15, left). The same trend is observed when ACT is applied in our dynamical network analysis for identifying interior critical residues (Supp. Fig. 15, right).

Thus, despite the different ways of modeling conformational change, we find that our conservation results hold for different models, thereby demonstrating that our method is general with respect to how motions are defined.

Critical Residues that Coincide with Human Disease Variants

Within our dataset of high-confidence alternative conformations, we identify 21 distinct proteins that are hit by known disease mutations, as collected from HGMD (Fig.

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: GN-

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: GN-

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: GN-

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: there are 176 distinct human proteins for which transcript IDs are available. Within this set of 176 distinct proteins,

4A) (Stenson et al 2014). Many of these proteins have been studied for their important biomedical significance. Examples include hemoglobin, phenylalanine hydroxylase, p53, and Ras (a full list of the affected PDBs, along with the afflicted residues, are provided as Supp. Files 2 and 3).

For 15 of this set of 21 proteins, the sites of HGMD mutations coincide with residues which lie in prioritized binding leverage surface residues. An example is Ras, shown in Fig. 4B). Likewise, 10 proteins have interior critical residues that overlap with sites of HGMD mutations. An example of such a system is p53, shown in Fig. 4C. The lists of structures for which prioritized binding sites or interior critical residues are affected by HGMD are provided in Supp. Files 4 and 5.

We note that, for several proteins, our identified critical residues coincide with known disease loci for which the mechanism of pathogenicity had been unclear (that is, some of our identified critical residues coincide with HGMD SNV locations, for which the pathogenic mechanism of the amino acid change associated with the HGMD is not obvious, but for which an allosteric mechanism provides a plausible means of understanding why the amino acid change may disrupt protein function). We emphasize that such disease loci constitute examples of “cryptic sites”, and that our framework helps to shed light on such regions in for which plausible alternative mechanisms of pathogenicity are not readily available.

A case-in-point is provided by a fibroblast growth factor receptor (Supp. Fig. 9), variants in which have been linked to Apert syndrome and Crouzon syndrome (diseases that manifest in craniofacial defects). Dotted lines in this plot highlight cryptic sites: loci that result in disease upon amino acid changes, for reasons that are not entirely clear (see Supp. Table 6). The incorporation of residues that fall in high binding-leverage sites or constitute high-betweenness loci in dynamic representations of the protein structure (i.e., our critical residues) adds an additional layer of annotation to the protein sequence, and these critical sites may help to explain poorly understood disease variants.

Finally, as we have done for HGMD SNVs, we also searched the NCBI ClinVar database (Landrum et al, 2014) for instances in which our identified critical residues coincide with disease locations. The affected proteins generally match those identified in

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: GN

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: GN

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: SNP

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: , but

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: which other biological annotations (such as post-translational modification sites, deeply buried regions, etc.) fail to coincide, and for which even the published literature introducing these disease mutations fails to propose a mechanism (see Supp.

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: HGMD

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: Other proteins and diseases for which poorly understood HGMD variants coincide with our identified critical residues include hemoglobin (alpha-thalassemia and erythrocytosis), protein deglycase (Parkinson disease and amyotrophic lateral sclerosis), phenylalanine hydroxylase (Phenylketonuria), UDP-glucose 4-epimerase (epimerase-deficiency galactosemia), isocitrate dehydrogenase (early-onset osteoarthritis), and HFE (haemochromatosis). A full listing, along with the results of manual literature curation of the associated disease variants, is provided in Supplementary Table 6.

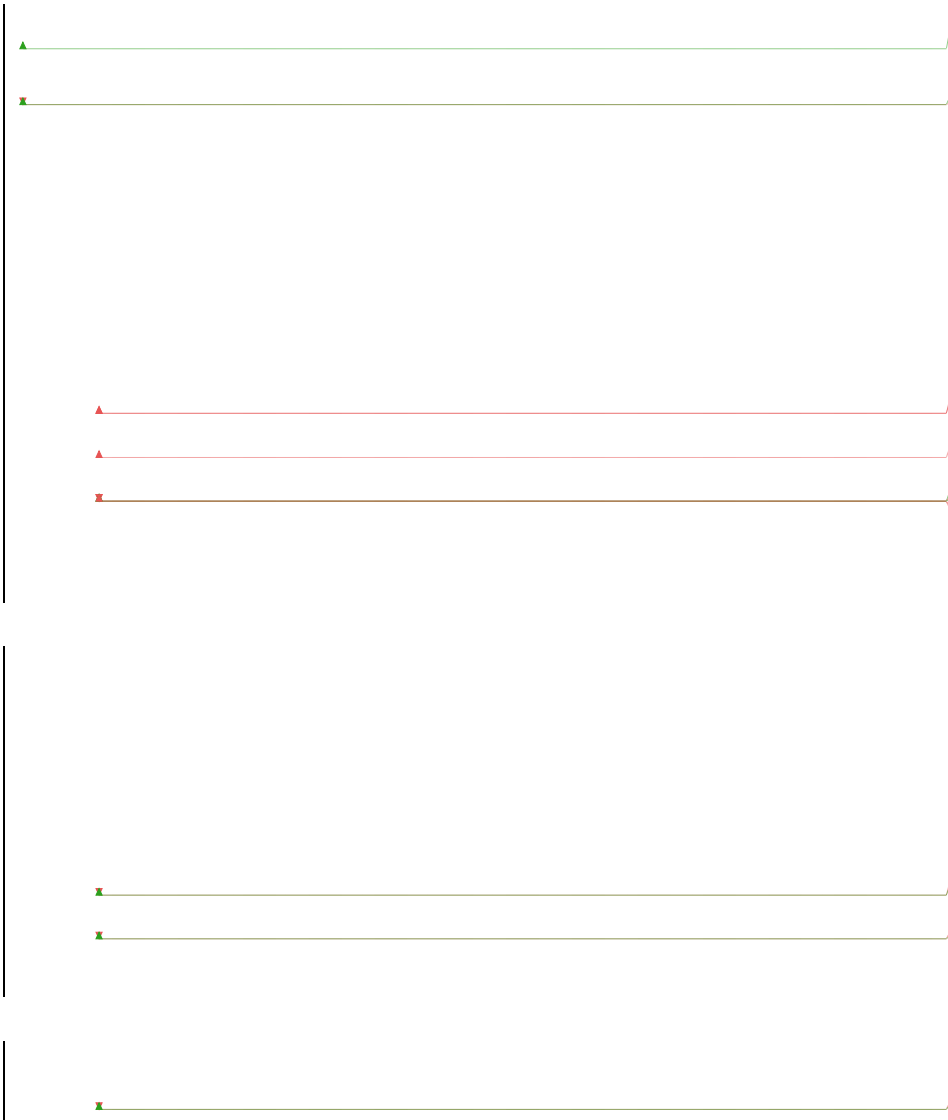
DECLAN CLARKE 7/19/15 12:42 AM
Deleted: SNPs

MSTH

the HGMD analysis above, and results are given in Supplementary Files 6 and 7 (for surface and interior critical residues, respectively).

Web Server (STRESS)

(Currently under development by Shantao and undergrad student Richard Chang).



DECLAN CLARKE 7/19/15 12:42 AM

Moved down [6]: METHODS
An overview of our pipeline is provided in Fig. 1, and we refer to this outline in the appropriate pipeline modules throughout. In brief, we perform MSAs for thousands of SCOP domains, with each alignment consisting of sequence-similar and sequence-identical domains. Within each alignment, we cluster the domains using structural similarity to determine the distinct conformational states.

DECLAN CLARKE 7/19/15 12:42 AM

Moved down [7]: Database-Wide Multiple Structure Alignments ... [14]

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: We then implement two complementary approaches for identif ... [13]

DECLAN CLARKE 7/19/15 12:42 AM

Formatted: Subscript

DECLAN CLARKE 7/19/15 12:42 AM

Formatted: Subscript

DECLAN CLARKE 7/19/15 12:42 AM

Formatted: Subscript

DECLAN CLARKE 7/19/15 12:42 AM

Moved down [8]: values for MSAs, as well as the motivating conceptual framewo ... [15]

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: values that exceed 1.

DECLAN CLARKE 7/19/15 12:42 AM

Moved down [9]: When performing K-means clustering with the gap statistic, ... [16]

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: The K

DECLAN CLARKE 7/19/15 12:42 AM

Formatted: Subscript

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: Once the optimal K value was determined for each of the N MSAs, w ... [17]

DECLAN CLARKE 7/19/15 12:42 AM

Moved down [10]: to generate dendrograms for each of the selected N ... [18]

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: The next step is to assign each unique domain to its respective cluster ... [19]

DECLAN CLARKE 7/19/15 12:42 AM

Moved down [11]: For each sequence group, we perform 1000 K-means clus ... [20]

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: et al

DECLAN CLARKE 7/19/15 12:42 AM

Moved down [12]: This method is motivated by the observation that allo ... [21]

DISCUSSION & CONCLUSIONS

The same principles of energy landscape theory that dictate protein folding have emerged as essential to understanding how folded proteins explore different conformational states in order to regulate needed functions. As in the case for folding, these landscapes are shaped not only by the protein sequence itself, but also by extrinsic conditions, such as post-translational modifications, interactions with ligands, substrates, or other proteins, or physiological conditions within the cell. Such external factors often regulate protein activity by introducing allosteric-induced changes, which ultimately reflect changes in the topology and population distributions of the energetic landscape.

In this regard, allostery provides a sensible platform from which to study protein behavior in the context of their energetic landscapes. Understanding allosteric signal transmission inevitably entails a consideration of the dynamic properties that generally accompany and are required for such allosteric behavior, as well as the identification of the residues that are essential for such behavior. Although not all allosteric proteins undergo conformational change (i.e., allosteric signals may sometimes be transmitted by changing the frequency of native-state fluctuations, rather than imparting large changes in conformational topology) [[cite Rodgers, Nussinov]] and not all conformational changes are associated with allostery [[cite ex: Calmodulin]], it is frequently the case that

- DECLAN CLARKE 7/19/15 12:42 AM
- Moved down [13]: for details reg... [23]
- DECLAN CLARKE 7/19/15 12:42 AM
- Deleted: We refer the reader to the v... [22]
- DECLAN CLARKE 7/19/15 12:42 AM
- Moved down [14]: was used to m... [24]
- DECLAN CLARKE 7/19/15 12:42 AM
- Deleted: A simple square well potential
- DECLAN CLARKE 7/19/15 12:42 AM
- Deleted: outlined by Mitternacht et ... [25]
- DECLAN CLARKE 7/19/15 12:42 AM
- Moved down [15]: When running ... [26]
- DECLAN CLARKE 7/19/15 12:42 AM
- Moved down [16]: an interaction ... [28]
- DECLAN CLARKE 7/19/15 12:42 AM
- Deleted: In the framework originally ... [27]
- DECLAN CLARKE 7/19/15 12:42 AM
- Deleted: For potential well widths, v... [29]
- DECLAN CLARKE 7/19/15 12:42 AM
- Moved down [17]: (attractive in t... [30]
- DECLAN CLARKE 7/19/15 12:42 AM
- Deleted: This benchmark set of prot... [31]
- DECLAN CLARKE 7/19/15 12:42 AM
- Moved down [18]: Using this app... [32]
- DECLAN CLARKE 7/19/15 12:42 AM
- Deleted: This weighting scheme is t... [33]
- DECLAN CLARKE 7/19/15 12:42 AM
- Moved down [19]: . This added ef... [34]
- DECLAN CLARKE 7/19/15 12:42 AM
- Moved down [20]: After weights ... [35]
- DECLAN CLARKE 7/19/15 12:42 AM
- Deleted: Residues that are involved ... [36]
- DECLAN CLARKE 7/19/15 12:42 AM
- Moved down [21]: Edge between ... [37]
- DECLAN CLARKE 7/19/15 12:42 AM
- Deleted: 1000 Genomes and HG... [38]
- DECLAN CLARKE 7/19/15 12:42 AM
- Moved down [22]: VCF files cont... [39]
- DECLAN CLARKE 7/19/15 12:42 AM
- Deleted: For nonsynonymous SNPs, ... [40]
- DECLAN CLARKE 7/19/15 12:42 AM
- Moved down [23]: on to protein ... [41]
- DECLAN CLARKE 7/19/15 12:42 AM
- Deleted: SNP to specific residues wi... [42]
- DECLAN CLARKE 7/19/15 12:42 AM
- Moved down [24]: As a quality as ... [43]
- DECLAN CLARKE 7/19/15 12:42 AM
- Deleted: SNPs in ExAC were analyz... [44]
- DECLAN CLARKE 7/19/15 12:42 AM
- Deleted: and...landscape theory that ... [45]
- DECLAN CLARKE 7/19/15 12:42 AM
- Deleted: principle ...esidues that are ... [46]

allosteric behavior is accompanied by substantial shifts in configurational space. Though a small number of examples in which allostery can occur without conformational change have been discussed in the literature (Tsai et al, 2009; Nussinov et al, 2015), the fact that these specific systems have been highlighted as exceptions underscores the important role played by conformational change in the vast majority of well-studied proteins, many of which have been investigated as a result of their significance in disease. In addition, we note that some proteins captured in our pipeline of alternative conformations do not exhibit allosteric behavior as part of their native functionality within cells, but the multiple energetic minima captured in their crystal structures may nevertheless be exploited for protein engineering [[cite C.J. Wilson, others]] or in pharmaceutical contexts [[cite]].

Molecular dynamics (MD) and NMR are some of the most common means of studying allostery and dynamic behavior. However, these methods have limitations when studying large and diverse protein datasets. Notably, MD is computationally very expensive, and is thus impractical when studying large numbers of proteins. Like MD, NMR yields important insights, but NMR structure determination is not only labor-intensive and best suited to specific classes of structures or dynamics (such as those with greater disordered content, or motions that operate on different time scales), but in addition, they constitute a relatively small fraction of the available structures (currently about 10%).

Given the limitations in applying MD, NMR, or related methods to large numbers of proteins, there remains a need to evaluate dynamic behavior in a systemized way across many proteins at once. This type of investigation applied to many proteins simultaneously also provides a means of better characterizing the large number of variants that have been shown to be deleterious through next-generation sequencing initiatives, thereby shedding light on the mechanisms at play for specific disease variants. Notably, such a database-scale approach is also much easier to exploit in studies focused on large networks of protein-protein interactions.

A database-scale approach necessitates careful and appropriate processing of data in the PDB. Though originally focused on finding structures for new proteins, there is now a great deal of redundancy in folds and proteins, and a concomitant greater degree of

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: structures of a particular nature,

heterogeneity from a functional point of view (i.e., more models for a given protein in different biological states). This redundancy opens the door to large-scale analyses aimed at investigating protein conformational heterogeneity and potential allosteric behavior on a database-level scale.

Thus, motivated by the idea that large differences in shape correspond to distinct conformations that occupy different energetic wells (Fig. 2), we describe and implement a pipeline for the identification of structures in distinct conformations using a statistical formalism. In doing so, we integrate data from the large number of X-ray crystal structures in the PDB, and simultaneously avoid the use of computationally expensive processes. The distinct conformations culled in this analysis are manually determined to correspond to proteins known to be in distinct functional states, such as active and inactive states, or holo and apo configurations. Users may submit protein structures to our server (STRESS) in order to perform their own analyses for identifying essential residues at the surface or within the interior. We emphasize that, as next-generation sequencing initiatives continue to provide a clearer picture of conservation at the residue level, structural biologists will increasingly find unexplained regions under strong selection. Our server readily enables the user to probe their own protein for potential allosteric regions, thereby helping to shed light on many of these regions.

Different conformations are used as the raw material for the identification of residues that may be important in the context of their allosteric behavior. We introduce a hybrid method to identify essential residues at the surface and the protein interior. To identify residues closer to the protein surface that may mediate allosteric behavior, we describe a modified version of the binding leverage method developed by Mitternacht and Berezovsky. We introduce information about the heavy atoms when searching the protein surface for sites in which the introduction of a ligand could strongly perturb the conformational changes of a protein, thereby finding sites that more closely reflect cavities in the protein topology. Secondly, after the sites are ranked by their ability to perturb the motions derived through anisotropic network models, we use a formalism originally used in the context of protein folding, the energy gap [[cite]], in order to define a threshold for selecting the high-confidence prioritized sites. We demonstrate that the set

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: ,

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: that, to our knowledge, has never previously been applied in the context of protein structures.

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: Though a database of allosteric proteins (Allosteric Database, ASD), has been described previously (Huang et al, 2011), the correspondence between ASD and our dataset exhibits relatively poor overlap. We emphasize that ASD was built using literature curation rather than direct considerations of the physical properties of the proteins themselves. In addition, the data in ASD is highly heterogeneous in nature, in that the structures vary considerably in terms of resolution and experimental origin, and paired entries need not be sequence-identical. Our focus is on a more confident set of large-scale conformational changes exclusively, with minimized noise. ... [47]

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: In the first of two complementary approaches for identifying such

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: (one approach being tailored to residues

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: , such as those in effector binding sites or allosteric sinks, and the other approach tailored to residues deeper inside

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: , which

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: transmit signals through distal regions),

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: et al.

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: candidate

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: candidate

of high-confidence sites overlaps reasonably well with known ligand binding sites for a set of well-studied canonical allosteric systems.

We employ a dynamical network-based analysis to search for sets of residues that may act as bottlenecks between communities in the protein structure. These communities are defined using the GN formalism, with edge weights reflecting the dynamic properties of the protein. This network-based analysis finds residues that are both internal to and within protein loops, some of which may be on the protein surface.

Thus, we emphasize that, while many previous studies use sequence characteristics or biophysical properties of individual amino acids to investigate allostery, focus on only the interior or surface residues, or may otherwise be restricted to a small number of proteins, we work on many proteins simultaneously within a generalized framework to use a *mechanistic* approach for identifying both surface and core residues which may be important for imparting allosteric behavior.

Our method is motivated to find many of the so-called cryptic elements in protein structures. Thus, we investigate the conservation of our critical residues in both inter-species and intra-human genomes contexts. The residues identified using the dynamical network analysis are shown to be conserved relative to other residues in the protein. More notably, the critical residues identified tend to be more conserved than residues with the same degree distribution (i.e., number of neighboring residues) within protein structures.

In addition, this greater conservation is also reflected in the genomes of modern-day humans: non-synonymous SNVs tend to hit these critical residues with lower frequencies than do other non-synonymous SNVs hitting the same protein, suggesting that amino acid changes at these critical sites may be more deleterious than changes in other parts of the protein. This trend is exhibited in both the distribution of allele frequencies, as well as with respect to the fraction of rare SNVs, and these results are observed when using either 1000 Genomes or ExAC datasets. We observe similar (though weaker) trends for surface residues implicated in allosteric behavior.

HGMD was used in order to identify any known disease-causing variants that hit the proteins in our dataset, and we found that several known disease SNVs, as culled from HGMD, hit the residues that we identify as being critical for allostery, on both the surface and within the interior. Given a particular set of PDBs (for example, the set of

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: In our second approach for finding allosteric residues, we

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: To evaluate the ability of this

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: identify residues that may be important for allosteric behavior

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: their

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: SNPs

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: SNPs

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: SNPs

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: SNPs

PDBs for which interior critical residues overlap with HGMD variant loci), the number of distinct proteins represented in this set was obtained by ensuring that no protein shares more than 90% sequence identity with any other protein in the set. We note that our set of 238 distinct HGMD proteins are those for which PDB structures are available, with the PDB structures satisfying structure quality criteria.

Given that allostery has previously been studied in the context of individual proteins, there are several notable implications of our database-scale analysis. That we achieved compelling results suggests that the level of coarse graining (i.e., in X-ray crystal structures and using ANMs instead of MD) was low enough to still recapitulate biologically interesting findings. That this pipeline can be applied en masse also suggests avenues for future applications, including applications to PPIs, guiding experimental studies to prioritize residues that are candidates for allosteric behavior (cite Rama Ranganathan, others), the simultaneous characterization of many disease variants in a diverse set of proteins (HGMD), and drug development pipelines/screens in which a drug is targeted to groups of functionally related proteins (such as those related to a particular signaling cascade or functional module) rather than to specific individual structures. Knowledge of predicted allosteric sites across many proteins may be used to identify the best proteins for which drugs should be engineered, as well as instances in which sequence variation is likely to have the greatest impact by modifying the relative populations of different states.

METHODS

An overview of our pipeline is provided in Fig. 1, and we refer to this outline in the appropriate pipeline modules throughout. In brief, we perform MSAs for thousands of SCOP domains, with each alignment consisting of sequence-similar and sequence-identical domains. Within each alignment, we cluster the domains using structural similarity to determine the distinct conformational states. We then implement coarse-grained models of protein motions to identify allosteric sites on the protein surface, as well as dynamical network analysis to identify allosteric residues internal to the protein.

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: GN

DECLAN CLARKE 7/19/15 12:42 AM

Moved (insertion) [6]

DECLAN CLARKE 7/19/15 12:42 AM

Moved (insertion) [7]

Database-Wide Multiple Structure Alignments

FASTA files of all SCOP domains were downloaded from the SCOP website (version 2.03) [cite]. In order to better ensure that large structural differences between sequence-identical or sequence-similar domains are a result of differing biological states (such as holo vs. apo, phosphorylated vs. unphosphorylated, etc.), and not an artifact of missing coordinates in X-ray crystal structures, the FASTA sequences used were those corresponding to the ATOM records of their respective PDBs. In total, this set comprises 162,517 FASTA sequences.

BLASTClust [cite] was downloaded from the NCBI database and used to organize these FASTA sequences into sequence-similar groups at seven levels of sequence identity (100%, 95%, 90%, 70%, 50%, 40%, and 30%). Thus, for instance, running BLASTClust with a parameter value of 100 provides a list of FASTA sequence groups such that each sequence within each group is 100% sequence identical, and in general, running BLASTClust with any given parameter value provides sequence groups such that each member within a group shares at least that specified degree of sequence identity with any other member of the same group (see top of Fig. 1).

To ensure that the X-Ray structures used in our downstream analysis are of sufficiently high quality, we removed all of those domains corresponding to PDB files with resolution values poorer than 2.8, as well as any PDB files with R-Free values poorer than 0.28. The question of how to set these quality thresholds is an important consideration, and was guided here by a combination of the thresholds conventionally used in other studies which rely on large datasets of structures [cite Kosloff 2008, Burra 2009, others], as well as the consideration that many interesting allosteric-related conformational changes may correlate with physical properties that sometimes render very high resolution values difficult (such as localized disorder or order-disorder transitions). As a result of applying these filters, 45,937 PDB IDs out of a total of 58,308 unique X-Ray structures (~79%) were kept for downstream analysis.

For each sequence-similar group at each of the seven levels of sequence identity, we performed multiple structure alignment (MSA) using only those domain structures that satisfy the criteria outlined above. Thus, the MSAs were generated only for those groups containing a minimum of two domains that pass the filtering criteria. The

STAMP and MultiSeq plugins of VMD were used to generate the MSAs. Heteroatoms were removed from each domain prior to performing the alignments.

The quality of the resultant MSA for each sequence-similar group depends on the root structure used in the alignment. To obtain the optimal MSA for each group of N domains, we generated N MSAs, with each alignment using a different one of the N domains as the root structure. The best MSA generated (as measured by STAMP's "sc" score) was taken as the MSA for that group. Note that, in order to aid in performing the MSAs, MultiSeq was used to generate sequence alignments for each group.

Finally, for each of the N MSAs generated, MultiSeq was used to calculate two measures of structural similarity between each pair of domains within a group: RMSD and Q_H . A fuller description of Q_H is provided in the Supplementary text.

For each group of sequence-similar domains, the final output of the structure alignment is a symmetric matrix representing all pairwise RMSD values (as well as a separate matrix representing all pairwise Q_H values) within that group. The matrices for all MSAs are then used as input to the K-means module.

The K -values for MSAs, as well as the motivating conceptual framework, are summarized in Fig 2. About 3000 different domains had a K -value of 1 (i.e., one conformation identified), whereas the K -values of close to 2000 domains exceed 1 (these exhibit multiple conformations, Fig. 2C). For proteins, close to 8000 had a K -value of 1, and about 1000 proteins had K -values that exceed 1. When performing K-means clustering with the gap statistic, very similar results were obtained when clustering structures on the basis of pairwise RMSD or pairwise Q_H (Supp. Fig. 3), so we use RMSD in our downstream analyses.

The fully-processed output for identifying high-confidence alternative conformations (which contains over 1100 proteins) is provided as a flat text file in the Supplementary content (Supp. File 1), and it is also included in our server as a downloadable text file. In addition to listing the alternative conformations, this file also lists descriptive statistics for each entry, including the RMSD between distinct conformers, cluster membership, degrees of confidence in assigning different structures to different clusters, etc.

DECLAN CLARKE 7/19/15 12:42 AM

Formatted: Subscript

DECLAN CLARKE 7/19/15 12:42 AM

Formatted: Subscript

DECLAN CLARKE 7/19/15 12:42 AM

Formatted: Subscript

DECLAN CLARKE 7/19/15 12:42 AM

Moved (insertion) [8]

DECLAN CLARKE 7/19/15 12:42 AM

Moved (insertion) [9]

DECLAN CLARKE 7/19/15 12:42 AM

Formatted: Subscript

We note that the pipeline above has been applied not only to SCOP domains, but also to individual proteins, with the only difference that only sequence-identical proteins were examined in this analysis.

Identifying Distinct Conformations in an Ensemble of Structures

For each MSA produced in the previous step, the corresponding matrix of pairwise RMSD values describes the degree and nature of structural heterogeneity among the crystal structures for a particular domain. The objective is to use this data in order to identify the biologically distinct conformations represented by an ensemble of structures. For a particular domain, there may be many available crystal structures. In total, these structures may actually represent only a small number of distinct biological states and conformations. For instance, there may be several crystal structures in which the domain is bound to its cognate ligand, while the remaining structures are in the apo state. Our framework for predicting the number of distinct conformational states in an ensemble of structures relies on a modified version of the K-means clustering algorithm.

A priori, performing K-means clustering assumes prior knowledge of the number of clusters (i.e., “K”) to describe a dataset. The purpose of K-means clustering with the gap statistic (Tibshirani et al, 2001) is to identify the optimal number of clusters intrinsic to a complex or noisy set of data points (which lie in N-dimensional space).

Given multiple resolved crystal structures for a given domain, this method (i.e., K-means with the gap statistic) estimates the number of conformational states represented in the ensemble of crystal structures (with these states presumably occupying different wells within the energetic landscape), thereby identifying proteins which are likely to undergo conformational change as part of their allosteric behavior.

As a first step toward clustering the structure ensemble represented by the RMSD matrix, it is necessary to convert this RMSD matrix (which explicitly represents only the relationships between distinct domains) into a form in which each domain is given its own set of coordinates. This step is necessary because the K-means algorithm acts directly on individual data points, rather than the distances between such points. Thus, we use multidimensional scaling [[ref Gower 1966 and Mardia, 1978]] to convert an N-by-N matrix (which provides all RMSD values between each pair of domains within a group of

N structures) into a set of N points, with each point representing a domain in (N-1)-dimensional space. The values of the N-1 coordinates assigned to each of these N points are such that the Euclidean distance between each pair of points are the same as the RMSD values in the original matrix. For an intuition into why N points must be mapped to (N-1)-dimensional space, consider an MSA between two structures. The RMSD between these two structures can be used to map the two domains to one-dimensional space, such that the distance between the points is the RMSD value. Similarly, an MSA of 3 domains may be mapped to 2-dimensional space in such a way that the pairwise distances are preserved; 4 domains may be mapped to 3-dimensional space, etc. The output of this multidimensional scaling is used as input to the K-means clustering with the gap statistic. We refer the reader to the work by Tibshirani et al for details governing how we perform K-means clustering with the gap statistic.

Once the optimal K-value was determined for each of the N MSAs, we confirmed that these values accurately reflect the number of clusters by manually studying several randomly-selected MSAs, as well as several MSAs corresponding of domain groups known to constitute distinct conformations (we also examined several negative controls, such as CAP, an allosteric protein which does not undergo conformational change [[ref]]).

To validate the output generated by this clustering algorithm, we manually annotated the alignments of a vast array well-studied canonical allosteric domains and proteins. There may be many factors driving conformational change, and those cases for which the change is induced by the binding to a simple ligand (i.e., a consideration of apo or holo states) constitute only a very small subset of the conformational shifts observed in the PDB. For instance, such shifts often result from protein-protein or protein-nucleic acid interactions, changes in oxidation states or in pH, mutations, binding to very large and complex ligands or the potential to bind to variable sets of ligands, post-translational modifications, interactions with the membrane, shifts in oligomerization states or configuration, etc. The gap statistic performed well in discriminating crystal structures that constitute such a diverse set, and this method has been validated using both domains (Supp. Figs. 4a-f) and protein chains (Supp. Figs. 4g-x).

RMSD values were used to generate dendrograms for each of the selected MSAs. The dendrograms are constructed using the hierarchical clustering algorithm built into R, hclust. [[ref Murtagh 1985]], with UPGMA (mean values) used as the chosen agglomeration method[[ref Sokal et al, 1958]].

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [10]

Each domain is assigned to its respective cluster using the assigned optimal K-values as input to Lloyd's algorithm. For each sequence group, we perform 1000 K-means clustering simulations on the MDS coordinates, and take the most common partition generated in these simulations to assign each protein to its respective cluster.

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [11]

We then select a representative domain from each of the assigned clusters. The representative member for each cluster is the member with the lowest Euclidean distance to the cluster mean, using the coordinates obtained by multidimensional scaling (see description above). These cluster representatives are then taken as the distinct conformations for this protein, and are used for the binding leverage calculations and networks analyses (below).

Modified Binding Leverage Framework

With the objective of identifying allosteric residues (specifically those on the protein surface), we employed a modified version of the binding leverage method for predicting likely ligand binding sites (Fig. 1, bottom-left), as described previously by Mitternacht and Berezovsky. This method is motivated by the observation that allosteric signals may be transmitted over large distances by a mechanism in which the allosteric ligand has a global affect on a protein's functionally important motions. For instance, introducing a bulky ligand into the site of an open pocket may disrupt large-scale motions if those motions normally entail that the pocket become completely collapsed in the apo protein. Such a modulation of the global motions may affect activity within sites that are distant from the allosteric ligand-binding site.

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [12]

We refer the reader to the work by Mitternacht and Berezovsky for details regarding the binding leverage method, though a general overview of the approach is given here. Hundreds or thousands of candidate allosteric sites are generated by simulations in which a simple ligand (comprising 2 to 8 atoms linked by bonds with fixed lengths but variable bond and dihedral angles) explores the protein's surface through

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [13]

many Monte Carlo steps. (Apo structures were used when probing protein surfaces for putative ligand binding sites). A simple square well potential (i.e., modeling hard-sphere interactions), was used to model the attractive and repulsive energy terms associated with the ligand's interaction with the surface. These energy terms depend only on the ligand atoms' distance to alpha carbon atoms in the protein, and they are blind to other heavy atoms or biophysical properties. Once these candidate sites have been produced, normal mode analysis is applied to generate a model of the apo protein's low-frequency motions. Each of the candidate sites is then scored based on the degree to which deformations in the site couple to the low-frequency modes; that is, those sites which are heavily deformed as a result of the normal mode fluctuations receive a high score (termed the binding leverage for that site), whereas sites which undergo minimal change over the course of a mode fluctuation receive a low binding leverage score. The list of candidate sites is then processed to remove redundancy, and then ranked based on this score. The model stipulates that the high-scoring sites are those that are more likely to be binding sites. Using knowledge of the experimentally-determined binding sites (i.e., from holo structures), the processed list of ranked sites is then used to evaluate predictive performance (see below).

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [14]

Our approach and set of applications differ from those previously developed in several key ways. When running Monte Carlo simulations to probe the protein surface and generate candidate binding sites, we used all heavy atoms in the protein when evaluating a ligand's affinity for each location. By including heavy atoms in this way (i.e., as oppose to using the protein's alpha carbon atoms exclusively), our hope is to generate a more realistic set of candidate ligand binding sites. Indeed, the exclusion of other heavy atoms leaves 'holes' in the protein which do not actually exist in the context of the dense topology of side chain atoms. Thus, by including all heavy atoms, we hope to reduce the number of false positive candidate sites, as well as more realistically model ligand binding affinities in general.

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [15]

In the framework originally outlined by Mitternacht and Berezovsky, an interaction between a ligand atom and an alpha carbon atom in the protein contributes -0.75 to the binding energy if the interaction distance is within the range of 5.5 to 8 Angstroms. Interaction distances greater than 8 Angstroms do not contribute to the

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [16]

binding energy, but distances in the range of 5.0 to 5.5 are repulsive, and those between 4.5 to 5.0 Angstroms are strongly repulsive (distances below 4.5 Angstroms are not permitted).

However, given the much higher density of atoms interacting with the ligand in our all-heavy atom model of each protein, it is necessary to accordingly change the energy parameters associated with the ligand's binding affinity. In particular, we varied both the ranges of favorable and unfavorable interactions, as well as the attractive and repulsive energies themselves (that is, we varied both the square well's width and depth when evaluating the ligand's affinity for a given site).

For well depths, we employed models using attractive potentials ranging from -0.05 to -0.75, including all intermediate factors of 0.05. For potential well widths, we tried performing the ligand simulations using the cutoff distances originally used by Mitternacht and Berezovsky (attractive in the range of 5.5 to 8.0 Angstroms, repulsive in the range of 5.0 to 5.5, and strongly repulsive in the range of 4.5 to 5.0). However, these cutoffs, which were originally devised to model the ligand's affinity to the alpha carbon atom skeleton alone, were observed to be inappropriate when including all heavy atoms. Thus, we also performed the simulations using a revised set of cutoffs, with attractive interactions in the range of 3.5 to 4.5 Angstroms, repulsive interactions in the range of 3.0 to 3.5 Angstroms, and strongly repulsive interactions in the range of 2.5 to 3.0 Angstroms.

In order to identify the optimal set of parameters for defining the potential function, we determined which combination of parameters best predicts the known binding sites for several well-annotated ligand-binding proteins. This benchmark set of proteins comprised threonine synthase (1E5X), phosphoribosyltransferase (IXTT), tyrosine phosphatase (2HNP), arginine kinase (3JU5), and adenylate kinase (4AKE). Using this approach, an attractive term of -0.35 for ligand-protein atom interactions within the range of 3.5 to 4.5 Angstroms was determined to be the best overall.

The biological assembly files (as well as individual proteins and standard PDBs) for several well-annotated allosteric and ligand-binding proteins [[list]] were downloaded from the Protein Data Bank (PDB). These proteins were chosen on the basis of literature

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [17]

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [18]

uration. Analyzed more proteins as gold standard (from several refs). Results are provided on server.

Network Analysis

In our implementation of the Girvan-Newman framework, edges between residues within a structure are drawn between any two residues that have at least one heavy atom within a distance of 4.5 Angstroms (excluding adjacent residues in sequence, which are not considered to be in contact). Network edges are weighted on the basis of their correlated motions, with the motions provided by anisotropic network models. We emphasize that, although the use of ANMs is more coarse-grained than MD, our use of ANMs is motivated by their much faster computational efficiency. This added efficiency is a required feature for our database-scale analysis.

Specifically, the weight w_{ij} between residues i and j is set to $-\log(|C_{ij}|)$. A high correlated motion between residues suggests strong information flow (see earlier discussion), and would thus result in a low value for w_{ij} . The 'distance' between residues i and j are thus taken to be very short, and this short distance means that any path involving this pair of residues is shorter as a result, thereby more likely placing this pair of residues within any given shortest path, and thus more likely rendering this pair of residues a bottleneck pair (thus, a high correlation results in a short distance, thereby more likely placing this pair in a short path which is thus more essential for intra-protein communication).

After weights are assigned, the betweenness for each edge is calculated. Residues that are involved in the highest-betweenness interactions connecting pairs of interacting communities are assigned to be in the class interior critical residues. Edge betweenness is defined as the total sum of shortest paths in which that edge is involved, with path lengths equal to the sum of edge weights (see Sethi et al, 2009 for a more detailed discussion).

Conservation Analyses

All cross-species conservation scores represent the ConSurf scores, as taken from the ConSurf Server [[cite]], in which scores for each protein chain are normalized to 0.

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [19]

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [20]

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [21]

Low (negative) ConSurf scores represent a stronger degree of conservation, and high (positive) scores designate less stringent selection.

All SNVs hitting protein-coding regions that result in amino acids changes (i.e., nonsynonymous SNVs) were collected from The 1000 Genomes Project (phase 3 release) [\[\[cite\]\]](#). VCF files containing the annotated variants were generated using VAT [\[\[cite\]\]](#). For nonsynonymous SNVs, the VCF files included the residue ID of the affected residue, as well as additional information (such as the corresponding allele frequency and residue type). To map the 1000 Genomes SNVs on to protein structures, FASTA files corresponding to the translated chain(s) of the respective transcript ID(s) were obtained using BioMart [\[\[cite\]\]](#). FASTA files for each of the PDB structures associated with these transcript IDs (the PDB ID-transcript ID correspondence was also obtained using BioMart) were generated based on the ATOM records of the PDB files. For each given protein chain, BLAST was used to align the FASTA file obtained from BioMart with that generated from the PDB structure. The residue-residue correspondence obtained from these alignments was then used in order to map each SNV to specific residues within the PDB. As a quality assurance mechanism, we confirmed that the residue type reported in the VCF file matched that specified in the PDB file.

ExAC variants were downloaded from the ExAC Browser (Beta), as hosted at the Broad Institute. Variants were mapped to all PDBs following the same protocol as that used to map 1000G variants, and only non-synonymous SNVs in ExAC were analyzed. When evaluating SNVs from the ExAC dataset, minor allele frequencies were used instead of DAF values (the ancestral allele is not provided in the ExAC dataset – thus, analysis is performed for MAF rather than DAF. However, we note that very little difference was observed when using AF or DAF values with 1000G data, and we believe that the results with MAF values would generally be the same to those with DAF values). Only structures for which at least one critical residue and one non-critical residue are hit by ExAC SNVs are included in the analysis (as with the 1000 Genomes analysis, this enables a more direct comparison between critical and non-critical residues, as comparisons between two different proteins would rely on the assumption of equal degrees of selection between such proteins).

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [22]

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [23]

DECLAN CLARKE 7/19/15 12:42 AM
Moved (insertion) [24]

FIGURE CAPTIONS

Figure 1

Pipeline for identifying distinct conformational states. *Top to bottom:* **a)** BLAST-CLUST is applied to the sequences corresponding to a filtered set of protein domains, thereby providing a large number of “sequence groups”, with each group being characterized by a high degree of sequence homology. **b)** For each sequence group, a multiple structure alignment of the domains is performed using STAMP (the example shown here is adenylate kinase. The SCOP IDs of the cyan domains, which constitute the holo structure, are d3hpqb1, d3hpqa1, d2eckb1, d2ecka1, d1akeb1, and d1akea1. The IDs of the apo domains, in red, are d4akea1 and d4akeb1). **c)** Using the pairwise RMSD values in this structure alignment, the structures are clustered using the UPGMA algorithm, K-means with the gap statistic (δ) is performed to identify the number of distinct conformations (2 in this example; more detailed descriptions of the graph are provided in the text and in Fig X). **d)** The domains which exhibit multiple structural clusters (i.e., those with a $\delta > X$ and $K > 1$) are then probed for the presence of strong allosteric sites, using binding leverage and dynamical network analysis (see Methods).

Figure 2

K-means clustering algorithm with the gap statistic. Number of binding sites per domain (a) and complex (b); c) An example dendrogram and respective structures of a multiple-structure alignment, with similarity measured by RMSD. The example shown is for phosphotransferase, and the K-means algorithm with the gap statistic identifies $K=2$ different conformational states (manually determined to represent the holo and apo states of phosphotransferase); **d)** Histograms representing the K_{δ} -values obtained across the database of SCOP domains and **e)** across PDB chains. Shown in (f) is a linear annotation diagram for fibroblast growth factor receptor. Shown is chain E of the PDB 1IIL, which corresponds to the FGFR2. Dotted lines highlight loci that correspond to HGMD sites that coincide with critical residues, but for which other annotations fail to coincide. Deeply-buried residues are defined to be those that exhibit a relative solvent-exposed surface area of 5% or less, and binding site residues are defined as those for which at

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: the complementary methods of

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: a) A schematized rendering

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: the k-means clustering algorithm;

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: c

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: k

DECLAN CLARKE 7/19/15 12:42 AM

Deleted: d

least one heavy atom falls within 4.5 Angstroms of any heavy atom in the binding partner (heparin-binding growth factor 2). The loci of PTM sites were taken from UniProt (accession no. P21802).

Figure 3

Conservation of predicted allosteric residues. **a)** Image of phosphofructokinase (PDB ID 3PFK), with red denoting sites with high binding leverage scores, and blue denoting sites with low scores. Known biological ligands are shown in white VDW rendering; **b)** Database-wide distributions of derived allele frequency (DAF) values of surface critical residues (red) and non-critical residues (blue); **c)** Corresponding distributions of conservation scores for surface critical residues (red) and non-critical residues (blue) **d)** Rendering of phosphofructokinase with interior critical residues highlighted as red spheres; **e)** Database-wide distributions of DAF values of interior critical residues (red) and non-critical residues (blue) **f)** Corresponding distributions of conservation scores for interior critical residues (red) and non-critical residues (blue).

Figure 4

HGMD Analyses. **a)** Venn diagram illustrating the number of distinct proteins in various categories; **b)** Ras (PDB ID 1NVV) is an example of a protein for which HGMD locations coincide with prioritized sites. Surface critical residues are shown as red spheres, and HGMD locations are in orange; **c)** p53 (PDB ID 2VUK) is an example of a protein for which HGMD locations coincide with interior critical residues. Interior critical residues that coincide with HGMD SNVs (red), critical residues that do not correspond with HGMD loci (green), and HGMD SNVs in non-critical residues (orange) are shown in VDW spheres.

REFERENCES

Arora, Karunesh, and Charles L. Brooks. "Large-scale allosteric conformational transitions of adenylate kinase appear to involve a population-shift mechanism." *Proceedings of the National Academy of Sciences* 104.47 (2007): 18496-18501.

Ashkenazy, Haim, et al. "ConSurf 2010: calculating evolutionary conservation in sequence and structure of proteins and nucleic acids." *Nucleic acids research* (2010): gkq399.

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: BL

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: ConSurf

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: BL

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: GN

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: GN

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: ConSurf

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: GN

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: BL

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: BL

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: GN

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: GN

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: SNPs

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: GN

DECLAN CLARKE 7/19/15 12:42 AM
Deleted: SNPs

Ashkenazy, Haim, Ron Unger, and Yossef Kliger. "Hidden conformations in protein structures." *Bioinformatics* 27.14 (2011): 1941-1947.

Bryngelson, Joseph D., et al. "Funnels, pathways, and the energy landscape of protein folding: a synthesis." *Proteins: Structure, Function, and Bioinformatics* 21.3 (1995): 167-195.

Bowman, Gregory R., et al. "Discovery of multiple hidden allosteric sites by combining Markov state models and experiments." *Proceedings of the National Academy of Sciences* 112.9 (2015): 2734-2739.

Burra, Prasad V., et al. "Global distribution of conformational states derived from redundant models in the PDB points to non-uniqueness of the protein structure." *Proceedings of the National Academy of Sciences* 106.26 (2009): 10505-10510.

Celniker, Gershon, et al. "ConSurf: using evolutionary data to raise testable hypotheses about protein function." *Israel Journal of Chemistry* 53.3 - 4 (2013): 199-206.

Dignam, John David, et al. "Allosteric interaction of nucleotides and tRNA^{Ala} with E. coli alanyl-tRNA synthetase." *Biochemistry* 50.45 (2011): 9886-9900.

Echols, Nathaniel, Duncan Milburn, and Mark Gerstein. "MolMovDB: analysis and visualization of conformational change and structural flexibility." *Nucleic Acids Research* 31.1 (2003): 478-482.

Exome Aggregation Consortium (ExAC), Cambridge, MA (URL: <http://exac.broadinstitute.org>) [May 2015]

Flicek P, Amode MR, Barrell D, Beal K, Brent S, et al. (2012) Ensembl 2012. *Nucleic Acids Res* 40: D84–90.

Flores, Samuel, et al. "The Database of Macromolecular Motions: new features added at the decade mark." *Nucleic acids research* 34.suppl 1 (2006): D296-D301.

Fox, Naomi K., Steven E. Brenner, and John-Marc Chandonia. "SCOPE: Structural Classification of Proteins—extended, integrating SCOP and ASTRAL data and classification of new structures." *Nucleic acids research* 42.D1 (2014): D304-D309.

Gerstein, Mark, and Werner Krebs. "A database of macromolecular motions." *Nucleic acids research* 26.18 (1998): 4280-4290.

Girvan, Michelle, and Mark EJ Newman. "Community structure in social and biological networks." *Proceedings of the National Academy of Sciences* 99.12 (2002): 7821-7826.

Glaser, Fabian, et al. "ConSurf: identification of functional regions in proteins by surface-mapping of phylogenetic information." *Bioinformatics* 19.1 (2003): 163-164.

Gunasekaran, K., Buyong Ma, and Ruth Nussinov. "Is allostery an intrinsic property of all dynamic proteins?" *Proteins: Structure, Function, and Bioinformatics* 57.3 (2004): 433-443.

Gower, J. C. (1966) Some distance properties of latent root and vector methods used in multivariate analysis. *Biometrika* 53, 325–328.

Grant, Gregory A., David J. Schuller, and Leonard J. Banaszak. "A model for the regulation of D - 3 - phosphoglycerate dehydrogenase, a Vmax - type allosteric enzyme." *Protein science* 5.1 (1996): 34-41.

Hubbard, Simon J., and Janet M. Thornton. "Naccess." Computer Program, Department of Biochemistry and Molecular Biology, University College London 2.1 (1993).

Huang, Zhimin, et al. "ASD: a comprehensive database of allosteric proteins and modulators." *Nucleic acids research* 39.suppl 1 (2011): D663-D669.

Kohl, Andreas, et al. "Allosteric inhibition of aminoglycoside phosphotransferase by a designed ankyrin repeat protein." *Structure* 13.8 (2005): 1131-1141

Kosloff, Mickey, and Rachel Kolodny. "Sequence - similar, structure - dissimilar protein pairs in the PDB." *Proteins: Structure, Function, and Bioinformatics* 71.2 (2008): 891-902.

Krebs, Werner G., and Mark Gerstein. "The morph server: a standardized system for analyzing and visualizing macromolecular motions in a database framework." *Nucleic Acids Research* 28.8 (2000): 1665-1675.

Lancichinetti, Andrea, and Santo Fortunato. "Community detection algorithms: a comparative analysis." *Physical review E* 80.5 (2009): 056117.

Landau, Meytal, et al. "ConSurf 2005: the projection of evolutionary conservation scores of residues on protein structures." *Nucleic acids research* 33.suppl 2 (2005): W299-W302.

Landrum MJ, Lee JM, Riley GR, Jang W, Rubinstein WS, Church DM, Maglott DR. ClinVar: public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Res.* 2014 Jan 1;42(1):D980-5. doi: 10.1093/nar/gkt1113. PubMed PMID: 24234437.

Laurent, M., et al. "Solution X-ray scattering studies of the yeast phosphofructokinase allosteric transition. Characterization of an ATP-induced conformation distinct in quaternary structure from the R and T states of the enzyme." *Journal of Biological Chemistry* 259.5 (1984): 3124-3126.

Liu, Ying, and Ivet Bahar. "Toward understanding allosteric signaling mechanisms in the ATPase domain of molecular chaperones." *Pacific Symposium on Biocomputing*. Vol. 15. 2010.

Manley, Gregory, Ivan Rivalta, and J. Patrick Loria. "Solution NMR and computational methods for understanding protein allostery." *The Journal of Physical Chemistry B* 117.11 (2013): 3063-3073.

Mardia, K.V. (1978) Some properties of classical multidimensional scaling. *Communications on Statistics – Theory and Methods*, A7, 1233–41.

Mitternacht, Simon, and Igor N. Berezovsky. "Binding leverage as a molecular basis for allosteric regulation." *PLoS computational biology* 7.9 (2011): e1002148.

Murtagh, F. (1985). "Multidimensional Clustering Algorithms", in *COMPSTAT Lectures 4*. Wuerzburg: Physica-Verlag (for algorithmic details of algorithms used).

Nussinov, Ruth, and Chung-Jung Tsai. "Allostery without a conformational change? Revisiting the paradigm." *Current opinion in structural biology* 30 (2015): 17-24.

- N Tibshirani, Robert, Guenther Walther, and Trevor Hastie. "Estimating the number of clusters in a data set via the gap statistic." *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 63.2 (2001): 411-423.
- Tsai, Chung-Jung, Buyong Ma, and Ruth Nussinov. "Folding and binding cascades: shifts in energy landscapes." *Proceedings of the National Academy of Sciences* 96.18 (1999): 9970-9972.
- Tsai, Chung-Jung, Antonio Del Sol, and Ruth Nussinov. "Allostery: absence of a change in shape does not imply that allostery is not at play." *Journal of molecular biology* 378.1 (2008): 1-11.
- Tsai, Chung-Jung, and Ruth Nussinov. "A unified view of "how allostery works"." (2014): e1003394.
- Rosvall, Martin, and Carl T. Bergstrom. "An information-theoretic framework for resolving community structure in complex networks." *Proceedings of the National Academy of Sciences* 104.18 (2007): 7327-7331.
- Sethi, Anurag, et al. "Dynamical networks in tRNA: protein complexes." *Proceedings of the National Academy of Sciences* 106.16 (2009): 6620-6625.
- Sethi, Anurag, et al. "A mechanistic understanding of allosteric immune escape pathways in the HIV-1 envelope glycoprotein." *PLoS computational biology* 9.5 (2013): e1003046.
- Sokal R and Michener C (1958). "A statistical method for evaluating systematic relationships". *University of Kansas Science Bulletin* 38: 1409-1438.
- Stenson et al (2014), The Human Gene Mutation Database: building a comprehensive mutation repository for clinical and molecular genetics, diagnostic testing and personalized genomic medicine. *Hum Genet* 133:1-9.
- [Watson, James D., and Francis HC Crick. "Molecular structure of nucleic acids." *Nature* 171.4356 \(1953\): 737-738.](#)
- Wiesmann, Christian, et al. "Allosteric inhibition of protein tyrosine phosphatase 1B." *Nature structural & molecular biology* 11.8 (2004): 730-737.
- Xiang, Yun, et al. "Simulating the effect of DNA polymerase mutations on transition-state energetics and fidelity: Evaluating amino acid group contribution and allosteric coupling for ionized residues in human pol β ." *Biochemistry* 45.23 (2006): 7036-7048.